



# Optimization of Positive Generalized Polynomials under $l^p$ Constraints

Laurent Baratchart, Marc Berthod, Loïc Pottier

## ► To cite this version:

Laurent Baratchart, Marc Berthod, Loïc Pottier. Optimization of Positive Generalized Polynomials under  $l^p$  Constraints. RR-2750, INRIA. 1995. inria-00073942

**HAL Id: inria-00073942**

**<https://inria.hal.science/inria-00073942>**

Submitted on 24 May 2006

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

INSTITUT NATIONAL DE RECHERCHE EN INFORMATIQUE ET EN AUTOMATIQUE

***Optimization of positive generalized  
polynomials under  $l^p$  constraints***

Laurent Baratchart, Marc Berthod, Loïc Pottier

**N° 2750**

Décembre 1995

PROGRAMME 4

 ***apport  
de recherche***



## Optimization of positive generalized polynomials under $l^p$ constraints

Laurent Baratchart\*, Marc Berthod\*\*, Loïc Pottier\*\*\*

Programme 4 — Robotique, image et vision  
Projet PASTIS

Rapport de recherche n° 2750 — Décembre 1995 — 38 pages

**Abstract:** The problem of maximizing a non-negative generalized polynomial of degree at most  $p$  on the  $l_p$ -sphere is shown to be equivalent to a concave one. Arguments where the *maximum* is attained are characterized in connection with the irreducible decomposition of the polynomial, and an application to the labelling problem is presented where these results are used to select the initial guess of a continuation method.

**Key-words:** constrained optimization, convex optimization, simulated annealing, relaxation, labeling

(Résumé : *tsvp*)

\*baratcha@sophia.inria.fr

\*\*berthod@sophia.inria.fr

\*\*\*pottier@sophia.inria.fr

# Optimisation de polynômes généralisés positifs sous contrainte de norme $l^p$

**Résumé :** On montre ici que la maximisation d'un polynôme généralisé non-négatif de degré  $p$  sur la sphère  $l - p$  est équivalent à un problème concave. La localisation du maximum est caractérisée en relation avec la décomposition du polynôme. Ces résultats sont appliqués au choix d'un point de départ pour une maximisation par une méthode de continuation, dans un problème d'optimisation combinatoire particulier: l'étiquetage

**Mots-clé :** optimisation sous contraintes, optimisation convexe, recuit simulé, relaxation, étiquetage

he problem of maximizing a non-negative generalized polynomial of degree at most  $p$  on the  $l_p$ -sphere is shown to be equivalent to a concave one. Arguments where the *maximum* is attained are characterized in connection with the irreducible decomposition of the polynomial, and an application to the labelling problem is presented where these results are used to select the initial guess of a continuation method.

## 1 Introduction

The main purpose of this paper is to analyse when the argument of the *maximum* on the  $l_p$ -sphere of a non-negative polynomial is unique, in the special case where the total degree of the polynomial does not exceed  $p$ . The author's original incentive to study this question lies with pattern recognition and image processing, where it turns out that maximizing a polynomial under suitable constraints is an effective way to approach certain combinatorial optimization problems that would hardly be tractable otherwise. This motivation is highlighted in section 2 where the classical labelling problem is sketched, as well as the *rationale* for replacing it by a continuous optimization scheme.

The issue just raised may be embodied into the more general problem of maximizing a non-negative generalized polynomial on the  $l_p$ -sphere. After some preliminaries on subhomogeneous functions and concavity in sections 3 and 4, we give in section 5 a systematic account of the solution when the degree is less than or equal to  $p$ . The approach is extremely elementary and consists in a simple change of variables that reduces the problem to concave maximization under linear constraints; determining strictly concave situations, however, involves a thorough discussion of the irreducible decomposition of polynomials which is linked to a non-linear eigenvalue problem of the Perron-Frobenius type. The results are partly carried over to  $l_p$ -constraints in product form in section 6.

We finally report in Section 7 on an application to finding the Maximum a Posteriori Mode in a Markov Random Field, for which numerical algorithms are also discussed briefly.

## 2 Some motivations: a deterministic approach to the labelling problem

The labelling problem arises quite naturally in pattern recognition and image processing. Actually, many image-interpretation tasks can be cast that way and the literature on this topic is plethoric. One of the first and best known instance of this phenomenon is the *Relaxation Labelling* approach that was celebrated in the computer vision community and for which we refer the reader to the seminal papers [10, 9] and [11]. A natural sequel to these developments was to rationalize the somewhat heuristic algorithms initially proposed, by defining merit functions to be optimized either locally [5] or globally [1]. An important step was taken when a well-founded probabilistic model was introduced by [7], which relied on the Hammersley-Clifford theorem [3]. As a matter of fact, all these references share a common framework that we illustrate here by describing the *Maximum a Posteriori Mode* problem (abbreviated as MAP) for a *Markov Random Field* (abbreviated as MRF).

We are given a set of units (or sites)  $\mathcal{S} = S_i$ ,  $1 \leq i \leq N$ , each of which may receive any label from 1 to  $M$ . A MRF on these units is defined as usual by a graph  $G$ , and the so-called clique potentials [7]. An edge of  $G$  connecting  $S_i$  and  $S_j$  is denoted by  $E_{ij}$  and  $V_i$  is the set of vertices (or sites) connected to a given vertex  $S_i$ . Let  $\mathcal{C}$  designate the set of all cliques of  $G$ , and define also  $C_i = \{c \in \mathcal{C}; S_i \in c\}$ . The number of sites in the clique  $c$  is its degree  $\deg(c)$ , and we set  $\deg(G) = \max_{c \in \mathcal{C}} \deg(c)$ .

A global discrete labelling  $L$  assigns one label  $L_i$  such that  $1 \leq L_i \leq M$  to each site  $S_i$  in  $\mathcal{S}$ . The restriction of  $L$  to the sites of a given clique  $c$  is denoted by  $L_c$ . The definition of the MRF is completed by the knowledge of the clique potentials  $V_{cL}$  (shorthand for  $V_{cL_c}$ ) for every  $c$  in  $\mathcal{C}$  and every  $L$  in  $\mathcal{L}$ , where  $\mathcal{L}$  is the set of the  $M^N$  discrete labelings.

In the MAP problem, the clique potentials stem from two sources of information: *a priori* knowledge about the restrictions that are imposed on the simultaneous labeling of connected neighboring units, and observations that were made on these units for a given occurrence of the problem. The goal is to find the labeling which maximizes the *a posteriori* probability given the observations.

Following Hammersley-Clifford, the probability of a given labeling  $L$  is given by:

$$P(L) \propto \prod_{c \in \mathcal{C}} \exp(-V_{cL}). \quad (1)$$

We assume here that the sufficient positivity condition for MRF is met *i.e.* that  $P(L) > 0$  for each  $L$ . It follows that solving the MAP problem amounts to find

$$\max_{L \in \mathcal{L}} \sum_{c \in \mathcal{C}} W_{cL}, \quad (2)$$

where  $W_{cL} = -V_{cL}$ .

*Deterministic Pseudo Annealing* (in short: DPA) has been proposed in [2] to tackle this maximization. The idea is to first replace the combinatorial question by an equivalent continuous optimization problem, and then try to solve this continuous problem by deforming it into a convex one.

More precisely, let us define  $f : \mathbf{R}^{NM} \rightarrow \mathbf{R}$  whose effect on

$$X = (x_{i,k})_{1 \leq i \leq N, 1 \leq k \leq M}$$

is given by

$$f(X) = \sum_{c \in \mathcal{C}} \sum_{l_c \in L_c} W_{cl_c} \prod_{j=1}^{\deg(c)} x_{c_j, l_{c_j}}, \quad (3)$$

where  $c_j$  denotes the  $j^{th}$  site of the clique  $c$  and  $l_{c_j}$  is the label assigned to it by  $l_c$ . It is clear from (3) that  $f$  is a polynomial in the  $NM$  variables  $x_{i,k}$ 's whose degree is  $\deg(G)$ . Moreover,  $f$  is linear in each variable separately. DPA in this case works as follows.

We define a compact subset  $\mathcal{K}$  of  $\mathbf{R}^{NM}$  by:

$$\forall i, k : x_{i,k} \geq 0 \quad \forall i : \sum_{k=1}^M x_{i,k} = 1.$$

The map  $f$  may have plenty of relative *maxima* on  $\mathcal{K}$ . However, there is always an absolute *maximum* attained on the boundary *i.e.* at some point  $X^*$  of the form:

$$\forall i, \exists k : x_{ik}^* = 1, l \neq k, \Rightarrow x_{il}^* = 0, \quad (4)$$



yielding naturally a discrete labelling. The difficulty is of course that standard search algorithms may typically lead to a local *maximum* and not to the absolute one. It is therefore of particular importance to find a good initial guess before applying the technique. This is precisely what DPA is designed for: we temporarily change the subset on which  $f$  is maximized so as to make the problem easy to solve, and then we track the *maximum* while *slowly* restoring the original constraints. At each step, the projection of the former point onto the new set of constraints is used as an initial guess for the next optimization.

To be specific, we trade  $\mathcal{K}$  for the set  $\mathcal{K}_p$  defined by

$$\forall i, k : x_{i,k} \geq 0 \quad \forall i : \sum_{k=1}^M x_{i,k}^p = 1.$$

There has been numerical evidence for a while that the *maximum* is attained at a single point when  $p > \deg(G)$ , and further that this point lies interior to  $\mathcal{K}_p$  [1]. The same holds true if  $\deg(G) = p$ , except for some degenerate zero-patterns of the coefficients for which the arguments of the *maximum* form a connected continuum intersecting certain coordinate axes. These facts are actually consequences of the results proved in section 5, making the case  $p \geq \deg(G)$  an easy problem. To achieve the DPA, it remains to decrease  $p$  down to 1, initializing the algorithm at each step from the projection onto the new set of constraints of the solution found at the preceding step. This is the heuristic part of the procedure, as we hope to track the right solution when bifurcations do occur.

As we now see, the labelling problem raises the issue of maximizing a non-negative polynomial on the positive face of a simplex, and our contribution to the DPA approach will consist in solving the analogous problem when the simplex is replaced by an  $l_p$ -sphere with  $p$  greater than or equal to the degree of the considered polynomial. The authors believe such a question possesses enough structure to make it worth studying, and other motivations like determining the dominant modes of certain nonlinear systems would also warrant such a study.

### 3 Preliminaries and notations.

We gather in this section a few pieces of notation and terminology which are of frequent use hereafter.

Let  $\mathbf{R}_+^n$  be the non-negative cone in  $\mathbf{R}^n$ , that is the subset of vectors with non-negative coordinates. The interior of  $\mathbf{R}_+^n$ , consisting of vectors with positive coordinates, will be denoted by  $\mathcal{P}_n$ . For  $x = (x_1 \cdots x_n) \in \mathbf{R}^n$  and  $p > 0$ , we denote by  $\|x\|_p = (\sum_i |x_i|^p)^{1/p}$  the  $l^p$  “norm” of  $x$ . This, actually, is a norm in the usual sense when  $p \geq 1$  only (otherwise triangular inequality will fail).

Given  $f : \Omega \rightarrow \mathbf{R}^m$ , where  $\Omega \subset \mathbf{R}^n$  is open, we denote by  $\partial f / \partial i(x)$  the partial derivative of  $f$  at  $x$  with respect to the  $i^{\text{th}}$  argument, and by  $\partial^2 f / \partial i \partial j(x)$  the second partial derivative with respect to the  $i^{\text{th}}$  and  $j^{\text{th}}$  arguments. If  $\mathcal{M}$  is a differentiable manifold and  $f : \mathcal{M} \rightarrow \mathbf{R}$  a differentiable function, we say that  $x \in \mathcal{M}$  is a *critical point* of  $f$  if the derivative  $Df(x)$  (which is defined on the tangent space  $\mathcal{T}_x \mathcal{M}$  to  $\mathcal{M}$  at  $x$ ) vanishes identically. When  $\mathcal{M}$  is embedded in  $\mathbf{R}^n$  and  $f$  extends to a differentiable function  $\tilde{f}$  in a neighborhood of  $x$  in  $\mathbf{R}^n$ , then  $x$  is critical if and only if the gradient vector of  $\tilde{f}$  at  $x$  is normal to  $\mathcal{T}_x \mathcal{M}$ .

A function  $f : \Omega \rightarrow \mathbf{R}$ , where  $\Omega \subset \mathbf{R}^n$  is open, is called *subhomogeneous of degree  $h \in \mathbf{R}$*  at  $x \in \Omega$  if there exists  $\epsilon(x) > 0$  such that

$$\forall \lambda \in [1, 1 + \epsilon(x)), \quad f(\lambda x) \leq \lambda^h f(x). \quad (5)$$

We say simply that  $f$  is subhomogeneous of degree  $h$  in  $\Omega$  if it is so at each point of  $\Omega$ . Here, a few comments are in order. Firstly, we restrict ourselves to  $\lambda \geq 1$  in the definition because if we allowed  $\lambda \in ]1 - \epsilon(x), 1 + \epsilon(x)[$ ,  $f$  would automatically be homogeneous. Secondly, the degree  $h$  in the definition is by no means unique: if  $f$ , for instance, is positive, any  $h' \geq h$  works.

Now, in the same manner than homogeneity translates into Euler’s identity, subhomogeneity translates into Euler’s inequality for differentiable functions: we say that a differentiable function  $f : \Omega \rightarrow \mathbf{R}$  satisfies *Euler’s inequality in degree  $h$*  at  $x = (x_1 \cdots x_n)$ , for some  $h \in \mathbf{R}$ , if

$$\sum_{i=1}^n \frac{\partial f}{\partial i}(x) x_i \leq h f(x). \quad (6)$$

If (6) holds at every  $x$ , we simply say that  $f$  satisfies Euler’s inequality. The link between subhomogeneity and Euler’s inequality is given in the following lemma.

**Lemma 1** *If  $\Omega$  is open in  $\mathbf{R}^n$ , and  $f : \Omega \rightarrow \mathbf{R}$  is a differentiable function which is subhomogeneous of degree  $h$  at  $x = (x_1 \cdots x_n)$ , then Euler's inequality in degree  $h$  holds for  $f$  at  $x$ . Conversely, if  $f$  satisfies Euler's inequality in degree  $h$  at every point  $\lambda x$  with  $1 \leq \lambda < 1 + \epsilon(x)$  for some  $\epsilon(x) > 0$ , then,  $f$  is subhomogeneous of degree  $h$  at  $x$ .*

*Proof:* suppose  $f$  is subhomogeneous, and put  $g_x(\lambda) = \lambda^h f(x) - f(\lambda x)$  for  $x \in \Omega$ . The function  $g_x : [1, 1 + \epsilon(x)) \rightarrow \mathbf{R}$  is non-negative, and vanishes at 1, hence  $g'_x(1) \geq 0$ . Expanding the derivative yields (6). Conversely, assume that  $f$  satisfies (6) at each  $\lambda x$  with  $1 \leq \lambda < \epsilon(x)$ . For such  $\lambda$ 's, we have

$$\frac{d g_x}{d \lambda} = h \lambda^{h-1} f(x) - \sum_{i=1}^n \frac{\partial f}{\partial i}(\lambda x) x_i \geq \frac{h}{\lambda} [\lambda^h f(x) - f(\lambda x)] = \frac{h g_x(\lambda)}{\lambda}.$$

This means  $\lambda g'_x \geq h g_x$ , or equivalently  $(g_x / \lambda^h)' \geq 0$ . Now, since the function  $g_x / \lambda^h$  vanishes at 1 and has non-negative derivative on  $[1, 1 + \epsilon(x))$ , it is non-negative there and so is  $g_x$ .  $\square$

For  $p \neq 0$ , we define throughout  $\varphi_p : \mathcal{P}_n \rightarrow \mathcal{P}_n$  by putting

$$\varphi_p(x_1, \dots, x_n) = (x_1^{\frac{1}{p}}, \dots, x_n^{\frac{1}{p}}). \quad (7)$$

Clearly,  $\varphi_p$  is a diffeomorphism.

Let  $A = (a_{i,j})$  be a real  $n \times n$  matrix. It is called *irreducible* if two distinct indices  $i$  and  $j$  can always be linked by a chain  $i = i_1, \dots, i_k = j$  in such a way that  $a_{i_\ell, i_{\ell+1}} \neq 0$ .

Let  $I = \{1 \cdots n\}$  be the set of indices, and  $e_i$  be the  $i^{\text{th}}$  vector of the canonical basis of  $\mathbf{R}^n$ . If  $J \subset I$ , we shall denote by  $E_J$  the subspace of  $\mathbf{R}^n$  spanned by the  $e_j$ 's for  $j \in J$ . Obviously,  $E_J$  consists of those vectors  $v = (v_1 \cdots v_n)$  with  $v_j = 0$  if  $j \notin J$ . We call  $E_J$  the *coordinate subspace* of  $\mathbf{R}^n$  associated with  $J$ . If  $I_1 \cdots, I_k$  is a partition of  $I$ , there is an orthogonal decomposition  $\mathbf{R}^n = \sum_j E_{I_j}$ , and we write accordingly  $v = \sum_j v_{I_j}$  for  $v \in \mathbf{R}^n$ .

It is a simple observation [6] that the irreducibility of a matrix  $A$  is equivalent to the non-existence of a non-trivial  $A$ -invariant coordinate subspace (that is, distinct from  $\{0\}$  and  $\mathbf{R}^n$  itself).

Suppose now that  $S = (s_{i,j})$  is a  $n \times n$  symmetric matrix. It is not difficult to check that the set  $I = \{1 \cdots n\}$  of indices can be partitioned into classes

$I_1 \cdots I_k$  such that the submatrices  $S_{I_\ell} = (s_{i,j})_{i \in I_\ell, j \in I_\ell}$  are irreducible and also  $s_{i,j} = 0$  if  $i$  and  $j$  belong to distinct  $I_\ell$ 's. This partition, which we call the *irreducible partition* of  $S$ , is well-defined since it corresponds to the decomposition of  $\mathbf{R}^n$  into minimal  $S$ -invariant coordinate subspaces. It follows from the definition that  $S_{I_\ell}$  is the matrix of the restriction of  $S$  to  $E_{I_\ell}$  when the latter is endowed with the canonical basis.

## 4 A concavity property

This section is instrumental for the remaining of the paper. The main result is that a  $C^2$  map  $f : \mathcal{P}_n \rightarrow \mathbf{R}$  whose first partial derivatives satisfy Euler's inequality in degree  $p-1$  for some  $p \neq 0$ , and whose second partial derivatives are non-negative, is such that  $f \circ \varphi_p$  is concave. This is essentially the content of Theorem 1 below. We begin with a computational lemma.

Let  $f : \mathcal{P}_n \rightarrow \mathbf{R}$  be  $C^2$  map, and put

$$\Phi_p = f \circ \varphi_p : \mathcal{P}_n \rightarrow \mathbf{R}$$

for  $p$  a nonzero real number. Denote the second derivative of  $\Phi_p$  at  $x$  by  $D^2\Phi_p(x)$ ; it is a bilinear form on  $\mathbf{R}^n$  that we identify with the  $n \times n$  symmetric matrix whose entry  $(i, j)$  is  $\partial^2\Phi_p/\partial i \partial j(x)$ . We also introduce another  $n \times n$  matrix  $M_{f,p}(x)$  whose entries at the point  $x = (x_1 \cdots x_n) \in \mathcal{P}_n$  are defined by the formulae

$$[M_{f,p}(x)]_{i,i} = x_i^{1-2p} \left[ x_i \frac{\partial^2 f}{\partial i \partial i}(x) - (p-1) \frac{\partial f}{\partial i}(x) \right], \quad (8)$$

$$[M_{f,p}(x)]_{i,j} = x_i^{1-2p} \left[ x_j \frac{\partial^2 f}{\partial i \partial j}(x) \right], \quad \text{for } i \neq j. \quad (9)$$

**Lemma 2** *With the above notations,  $D^2\Phi_p(x)$  is conjugate to  $M_{f,p}(\varphi_p(x))/p^2$  at any point  $x \in \mathcal{P}_n$ . More precisely, we have:*

$$B^{-1}(x) D^2\Phi_p(x) B(x) = \frac{M_{f,p}(\varphi_p(x))}{p^2}, \quad (10)$$

where  $B(x)$  is the diagonal matrix  $\text{diag}\{x_i\}$ .

*Proof:* this is a simple computation. First, we write

$$\frac{\partial \Phi_p}{\partial i}(x) = \frac{\partial f}{\partial i}(\varphi_p(x)) \frac{1}{p} x_i^{\frac{1}{p}-1}, \quad (11)$$

whence

$$\frac{\partial^2 \Phi_p}{\partial i \partial i}(x) = \frac{\partial^2 f}{\partial i \partial i}(\varphi_p(x)) \left( \frac{1}{p} x_i^{\frac{1}{p}-1} \right)^2 + \frac{\partial f}{\partial i}(\varphi_p(x)) \frac{1}{p} \left( \frac{1}{p} - 1 \right) x_i^{\frac{1}{p}-2}. \quad (12)$$

This can be rearranged as

$$\frac{1}{p^2} x_i^{\frac{1}{p}-2} \left[ \frac{\partial^2 f}{\partial i \partial i}(\varphi_p(x)) x_i^{\frac{1}{p}} - (p-1) \frac{\partial f}{\partial i}(\varphi_p(x)) \right]. \quad (13)$$

If  $i \neq j$ , we get similarly

$$\frac{\partial^2 \Phi_p}{\partial i \partial j}(x) = \frac{1}{p^2} x_i^{\frac{1}{p}-1} \frac{\partial^2 f}{\partial i \partial j}(\varphi_p(x)) x_j^{\frac{1}{p}-1}. \quad (14)$$

Now, compute  $B^{-1}(x) D^2 \Phi_p(x) B(x)$ . Since the  $i^{\text{th}}$  row of  $D^2 \Phi_p(x)$  gets divided by  $x_i$  while the  $j^{\text{th}}$  column gets multiplied by  $x_j$ , the result is  $M_{f,p}(\varphi_p(x))/p^2$ .  $\square$

We are now in position to state:

**Theorem 1** *Let  $f : \mathcal{P}_n \rightarrow \mathbf{R}$  be a  $C^2$  map. For  $p \neq 0$ , define  $\varphi_p$  and  $\Phi_p$  as before, and let  $x = (x_1 \cdots x_n)$  be a point in  $\mathcal{P}_n$  such that each partial derivative  $\partial f / \partial i$  satisfies Euler's inequality in degree  $p-1$  at  $\varphi_p(x)$ , while each second partial derivative  $\partial^2 f / \partial i \partial j$  is non-negative at  $\varphi_p(x)$ . Let finally  $I_1 \cdots I_k$  denote the irreducible partition of  $D^2 \Phi_p(x)$ . Then  $D^2 \Phi_p(x)$  defines a non-positive quadratic form. It is negative definite unless there exists an  $\ell$  such that Euler's inequality for  $\partial f / \partial i$  at  $\varphi_p(x)$  is in fact an equality for every  $i \in I_\ell$ . If we let  $I' \subset \{1 \cdots k\}$  be the set of such  $\ell$ 's, the kernel of  $D^2 \Phi_p(x)$  is the subspace spanned by the  $x_{I_j}$ 's for  $j \in I'$ .*

*Proof:* from Lemma 2, we see that the eigenvectors of  $D^2 \Phi_p(x)$  are the images under  $B(x)$  of those of  $M_{f,p}(\varphi_p(x))/p^2$ , and that the eigenvalues of the two matrices differ only by a factor  $1/p^2$ . In particular, since  $D^2 \Phi_p(x)$  is

a symmetric matrix hence has real eigenvalues, so does  $M_{f,p}(\varphi_p(x))$ . Let us denote by  $(m_{i,j})$  the entries of  $M_{f,p}(\varphi_p(x))$ . By assumption, we have Euler's inequality for  $\partial f / \partial i$  at  $\varphi_p(x)$ :

$$\sum_{j=1}^n \frac{\partial^2 f}{\partial i \partial j}(\varphi_p(x)) x_j^{\frac{1}{p}} \leq (p-1) \frac{\partial f}{\partial i}(\varphi_p(x)). \quad (15)$$

Upon multiplying by the positive quantity  $x_i^{\frac{1}{p}-2}$ , we get by the very definition of  $(m_{i,j})$ :

$$\sum_{j=1}^n m_{i,j} \leq 0, \quad \forall i \in \{1 \cdots n\}. \quad (16)$$

By hypothesis, all partial second derivatives of  $f$  are non-negative at  $\varphi_p(x)$ , so we see from the definition that  $m_{i,j} \geq 0$  if  $i \neq j$ . Now, a well-known theorem of Gerschgorin (see e.g. [6]) tells us that every eigenvalue of  $M_{f,p}(\varphi_p(x))$  belongs to a disc centered at some  $m_{i,i}$  of radius  $\sum_{j \neq i} |m_{i,j}| = \sum_{j \neq i} m_{i,j}$ . By (16), all these eigenvalues lie in the left half-plane, and since they are real, they are non-positive. This shows that the eigenvalues of  $D^2\Phi_p(x)$  are also non-positive, and so is the associated quadratic form.

We now compute the kernel of  $D^2\Phi_p(x)$ . Since the latter is symmetric and  $E_{I_j}$  is  $D^2\Phi_p(x)$ -stable by definition, we first observe that  $(D^2\Phi_p(x)v)_{I_j} = D^2\Phi_p(x)v_{I_j}$  for any  $v \in \mathbf{R}^n$ . Therefore,  $v$  belongs to the kernel of  $D^2\Phi_p(x)$  if and only if  $v_{I_j}$  belongs to the kernel of the restriction of  $D^2\Phi_p(x)$  to  $E_{I_j}$  for all  $j \in \{1, \dots, k\}$ . Hence, it is enough to prove that the kernel of  $D^2\Phi_p(x)_{I_j}$ , if non-trivial, is generated by  $x_{I_j}$ , and that Euler's inequality, when applied to  $\partial f / \partial i$ , is then an equality for every  $i \in I_j$ . Let  $n_j$  denote the cardinality of  $I_j$ . Because of the relationship between  $D^2\Phi_p(x)$  and  $M_{f,p}(\varphi_p(x))$  asserted in Lemma 2, it is equivalent to show that the kernel of  $(M_{f,p}(\varphi_p(x)))_{I_j}$ , if nontrivial, is generated by the vector  $(1 \cdots 1)$  of size  $n_j$ , and still Euler's inequality is an equality for  $i \in I_j$ .

Suppose now that the kernel of  $(M_{f,p}(\varphi_p(x)))_{I_j}$  does not consist of zero alone, and let  $w = (w_1 \cdots w_{n_j})$  be a non-zero vector in this kernel. Choose and index  $i_0 \in I_j$  such that  $w_{i_0}$  is of maximum modulus. Replacing  $w$  by  $-w$  if necessary, we may assume that  $w_{i_0} > 0$ . Define  $K \subset I_j$  to be the subset of indices  $i$  such that  $w_i = w_{i_0}$ . Pick any  $i_1 \in K$ . If we write that the  $i_1^{th}$

component of  $(M_{f,p}(\varphi_p(x)))_{I_j}$   $w$  is zero, we get:

$$\sum_{i \in I_j} m_{i_1, i} w_i = 0. \quad (17)$$

Upon multiplying (16) by  $w_{i_1}$ , and taking into account the fact that  $m_{i_1, i} = 0$  if  $i \notin I_j$  (since the same holds true for the matrix  $D^2\Phi_p(x)$ ), we also have that

$$\sum_{i \in I_j} m_{i_1, i} w_{i_1} \leq 0. \quad (18)$$

Subtracting (17) from (18) yields

$$\sum_{i \in I_j, i \neq i_1} m_{i_1, i} (w_{i_1} - w_i) \leq 0. \quad (19)$$

But each term in the sum is non-negative, so they all vanish, and equality holds in (19), hence in (18), and in Euler's inequality for  $\partial f / \partial i_1$  as well. If  $i \notin K$ , then  $w_{i_1} > w_i$  in (19) and therefore  $m_{i_1, i} = 0$  whence also  $m_{i, i_1} = 0$  (because they are proportional). Since  $i_1$  was arbitrary in  $K$ , it follows that  $E_K \subset E_j$  is invariant under  $(M_{f,p}(\varphi_p(x)))_{I_j}$ . But this matrix is irreducible since  $D^2\Phi_p(x)_{I_j}$  is, by definition of  $I_j$ . Thus, we have  $K = I_j$ , and  $w_i = w_{i_0}$  for every  $i \in I_j$ . This is precisely what we wanted to show. Conversely, it is clear that if Euler's inequality is an equality for every  $i \in I_j$ , the vector  $(1 \cdots 1)$  clearly lies in the kernel of  $(M_{f,p}(\varphi_p(x)))_{I_j}$ .  $\square$

To recap, Theorem 1 asserts that a  $C^2$  function  $f : \mathcal{P}_n \rightarrow \mathbf{R}$  is such that  $f \circ \varphi_p$  is concave as soon as

- (i) the second partial derivatives are non-negative at every  $x \in \mathcal{P}_n$ ; this is equivalent to saying that the gradient of  $f$  is non-decreasing for the usual partial ordering of  $\mathcal{P}_n$ , i.e.  $x \geq y$  iff  $y - x \in \mathcal{P}_n$ .
- (ii) the partial derivatives satisfy Euler's inequality in degree  $p - 1$  on  $\mathcal{P}_n$ ; by Lemma 1, this is equivalent to the seemingly more natural property that the derivatives of  $f$  are subhomogeneous of degree  $p - 1$  on  $\mathcal{P}_n$ .

## 5 $l_p$ -constrained maximization of positive generalized polynomials of degree at most $p$

In this section, we apply the preceding results to the problem of maximizing a generalized polynomial with non-negative coefficients (see the definition below) on the  $l^p$ -sphere when  $p$  is not less than the degree. This allows us to describe uniqueness and positivity properties of the solution. The approach is completely elementary and simply consists in composing the polynomial with  $\varphi_p$  so as to be back to standard optimization of a concave function under linear constraints. Since we want to give complete answers on uniqueness, however, we need to analyse cases when strict concavity prevails and this requires a slightly lengthier discussion of the irreducible decomposition which makes this section somewhat reminiscent of the Perron-Frobenius theory for nonnegative matrices. In effect, at the end of the section, we use the critical point equation to derive some kind of nonlinear generalization for symmetric matrices of the Perron-Frobenius theorem.

Strictly speaking, the facts that we shall use about non-negative generalized polynomials are subhomogeneity of the derivatives and nonnegativity of the second derivatives when the exponents involved are not less than 1. *In fact, the results and the proofs can be adapted to any function sharing these properties.* In particular, everything in this section extends to *infinite* sums  $\sum c_\alpha x^\alpha$  where  $\alpha \in \mathbf{R}_+^n$  is bounded by  $p$  in  $l^1$ -norm and the coefficients  $c_\alpha \in \mathbf{R}_+$  decrease fast enough to ensure that the series converges absolutely for, say,  $\|x\|_p < 1 + \epsilon$  for some positive  $\epsilon$ . Nevertheless, we shall stick to the case of generalized polynomials an application of which was described in the introduction.

A *generalized polynomial* is a function  $P : \mathcal{P}_n \rightarrow \mathbf{R}$  of the form:

$$P(x) = \sum_{\alpha} c_{\alpha} x^{\alpha}, \quad (20)$$

where  $\alpha$  ranges over a finite set of  $\mathbf{R}_+^n$ , and  $x^\alpha$  stands for  $x_1^{\alpha_1} \cdots x_n^{\alpha_n}$  with  $\alpha = (\alpha_1 \cdots \alpha_n)$ . By definition, the degree of  $P$  is  $h = \max_{\alpha} h_{\alpha}$ , where  $h_{\alpha} = \sum_i \alpha_i$ . If  $h_{\alpha} = h, \forall \alpha$ , we call  $P$  an homogeneous generalized polynomial of degree  $h$ . By convention, the zero polynomial is homogeneous of any



degree. We say that  $P$  has non-negative coefficients if  $c_\alpha \geq 0$  for all  $\alpha$ .

Now, we need to extend the notion of irreducibility, already introduced for matrices, to generalized polynomials. If  $P = \sum c_\alpha x^\alpha$  is such a polynomial, we can associate to  $P$  a graph whose vertices are the variables, and an edge connects two variables  $x_i$  and  $x_j$  iff there exists a term  $c_\alpha \neq 0$  with  $\alpha_i \neq 0$  and  $\alpha_j \neq 0$ . The adjacency matrix of this graph is symmetric and its irreducible partition is also called the irreducible partition of  $P$ . Another way to look at things is to observe that any generalized polynomial  $P$  in  $n$  variables can be written (in possibly many ways) as

$$P(x) = \sum_{j=1}^k P_{I_j}(x_{I_j}), \quad (21)$$

where the  $I_j$ 's, for  $j \in \{1 \cdots k\}$ , partition the set  $I = \{1 \cdots n\}$ ; it is easy to check that there is a unique minimal such decomposition (i.e. one that cannot be refined), where the  $P_{I_j}$ 's are defined up to a constant term and where the  $I_j$ 's are nothing but the irreducible partition of  $P$  already defined. If, in addition,  $P$  has nonnegative coefficients, this irreducible partition is that of the second derivative  $D^2P$  at any (and thus every) point of  $\mathcal{P}_n$ . The additive decomposition (21) associated to the irreducible partition is called the *irreducible decomposition* of  $P$ , and the polynomials  $P_{I_j}$  in this decomposition are called the *irreducible components*. Such a component may well be zero; in this case, the polynomial depends on fewer variables. Note also that irreducible components are defined only up to constant terms so that *any qualification concerning them should be understood modulo a constant term*. An irreducible component which is not  $P$  itself is said to be proper, and we say that  $P$  is irreducible if and only if it has no proper irreducible component. Equivalently, this means that the irreducible partition has only one element, namely  $I$  itself.

A family of functions to which theorem 1 applies naturally is the family of generalized polynomials with non-negative coefficients; this is due to the following result.

**Proposition 1** *A generalized polynomial  $P$  of degree  $h$  in  $n$  variables with non-negative coefficients is subhomogeneous of degree  $h$  in the positive cone*

$\mathcal{P}_n$ . If  $P$  is not homogeneous, it is in fact subhomogeneous of degree strictly less than  $h$  at any point of  $\mathcal{P}_n$ .

*Proof:* write  $P = \sum_{\beta} P_{\beta}$ , where each  $P_{\beta}$  is homogeneous of degree  $\beta$ . For any  $x \in \mathbf{R}^n$  and  $\lambda \in \mathbf{R}$ , we have

$$P(\lambda x) = \sum_{\beta} \lambda^{\beta} P_{\beta}(x). \quad (22)$$

Now for  $x \in \mathcal{P}_n$ , each  $P_{\beta}(x)$  is obviously non-negative, and it is clear if  $\lambda \geq 1$  that  $\lambda^{\beta} \leq \lambda^h$ ; therefore,  $P(\lambda x) \leq \lambda^h P(x)$ , showing that  $P$  is subhomogeneous of degree  $h$  in  $\mathcal{P}_n$ .

Assume now that  $P$  is not homogeneous, and let  $\beta_0$  be largest among those  $\beta$ 's such that  $P_{\beta} \neq 0$  and  $\beta < h$ . Pick  $\mu_0 > 0$  so small that  $h - \beta_0 - \mu > 0$ ; for  $x \in \mathcal{P}_n$ ,  $\lambda \geq 1$  and  $0 < \mu < \mu_0$ , put

$$g_{x,\mu}(\lambda) = \lambda^{h-\mu} P(x) - P(\lambda x). \quad (23)$$

We compute

$$g_{x,\mu}(\lambda) = \sum_{\beta} [\lambda^{h-\mu} - \lambda^{\beta}] P_{\beta}(x), \quad (24)$$

so that

$$\frac{dg_{x,\mu}}{d\lambda}(1) = \sum_{\beta} [h - \beta - \mu] P_{\beta}(x) \quad (25)$$

which is, by the definition of  $\mu$ , bounded from below by

$$(h - \beta_0 - \mu) P_{\beta_0}(x) - \mu P_h(x). \quad (26)$$

Since each component of  $x$  is positive, it follows that  $P_{\beta_0}(x) > 0$ , hence the above expression can be made positive by choosing  $\mu$  sufficiently small. Then  $g_{x,\mu}(\lambda)$  vanishes with positive derivative at  $\lambda = 1$ . Therefore, it remains non-negative on some interval  $[1, 1 + \epsilon(x))$ , so that  $P$  is subhomogeneous of degree  $h - \mu < h$  at  $x$ .  $\square$

As an application of Proposition 1 and Theorem 1, we study for later use the concavity properties of non-negative generalized polynomials when the degree does not exceed 1.

**Theorem 2** *A generalized polynomial  $P$  with non negative coefficients of degree at most 1 is concave on  $\mathbf{R}_+^n$ . If  $P(x) = \sum_{\ell=1}^k P_{I_\ell}$  is the irreducible decomposition, the kernel of  $D^2P(x)$  at  $x \in \mathcal{P}_n$  is the linear span of those  $x_{I_\ell}$ 's such that  $P_{I_\ell}$  is homogeneous of degree 1. In particular,  $P$  is strictly concave on  $\mathcal{P}_n$  if, and only if, it has no irreducible component which is homogeneous of degree 1, in which case the second derivative is negative definite at each point of  $\mathcal{P}_n$ .*

*Proof:* Since  $P$  is continuous on  $\mathbf{R}_+^n$  (for the exponents are non-negative), it is enough to show that  $P$  is concave on  $\mathcal{P}_n$ . Let  $\alpha_m$  be the smallest non-zero exponent for a variable appearing in  $P$ . Set  $p = \alpha_m^{-1}$  and define  $f = P \circ \varphi_p^{-1} = P \circ \varphi_{1/p}$ , so that  $f$  is obtained from  $P$  by changing each variable into its  $p^{th}$  power. Then,  $f$  is again a generalized polynomial with non-negative coefficients, of degree  $h \leq p$  (this is where we use  $\deg P \leq 1$ ), each variable of which appears at every occurrence with exponent at least 1. Thus, every partial derivative  $\partial f / \partial i$  of  $f$  is a generalized polynomial with non-negative coefficients, of degree at most  $h - 1$ . By Proposition 1 and Lemma 1,  $\partial f / \partial i$  satisfies Euler's inequality in degree  $h - 1$  and so *a fortiori* in degree  $p - 1$  on  $\mathcal{P}_n$ . Since the second partial derivatives of  $f$  are clearly non-negative on  $\mathcal{P}_n$ , Theorem 1 implies that  $P = f \circ \varphi_p$  has a non-positive second derivative at every point of  $\mathcal{P}_n$ , hence is concave there.

Let us now determine the kernel of the second derivative. Since  $P$  is non-negative, its irreducible partition is also that of  $D^2P(x)$  at any  $x \in \mathcal{P}_n$  and we deduce from Theorem 1 again that the kernel of this matrix is generated by those  $x_{I_\ell}$ 's such that  $\partial f / \partial i$  satisfies Euler's *equality* in degree  $p - 1$  at  $\varphi_p(x)$  for each  $i \in I_\ell$ . By construction, the irreducible partitions of  $P$  and  $f$  are identical, so we can write

$$f = \sum_{\ell=1}^k f_{I_\ell},$$

and it is clear that  $f_{I_\ell}$  is homogeneous of degree  $p$  if, and only if,  $P_{I_\ell}$  is homogeneous of degree 1. Hence, it remains for us to show that  $\partial f_{I_\ell} / \partial i$  satisfies Euler's identity in degree  $p - 1$  at  $\varphi_p(x)$  for each  $i \in I_\ell$  if, and only if,  $f_{I_\ell}$  is homogeneous of degree  $p$ . Sufficiency is obvious. By Proposition 1, necessity amounts to prove that a generalized polynomial which is not homogeneous of degree  $p$  (up to a constant term) cannot have partial derivatives each of which is homogeneous of degree  $p - 1$ . This is easy: arguing by contradiction, sup-

pose that  $Q(y_1, \dots, y_s)$  is such a polynomial and consider the homogeneous generalized polynomial  $\partial Q / \partial y_1$  of degree  $p - 1$ . By elementary integration, we get

$$Q(y_1, \dots, y_s) = Q_1(y_1, \dots, y_s) + Q_2(y_2, \dots, y_s),$$

where  $Q_1$  is homogeneous of degree  $p$  and  $Q_2$  does not depend on  $y_1$ . Now, for  $2 \leq i \leq s$ ,

$$\frac{\partial Q_2}{\partial y_i} = \frac{\partial Q}{\partial y_i} - \frac{\partial Q_1}{\partial y_i}$$

is homogeneous of degree  $p - 1$ , so we can iterate the process and write  $Q_2 = Q_1^1 + Q_2^1$  where  $Q_1^1$  is homogeneous of degree  $p$ , and  $Q_2^1$  depends on  $y_3, \dots, y_s$  only, while still having homogeneous derivatives of degree  $p - 1$ . By induction, we conclude that

$$Q = Q_1 + Q_1^1 + \dots + Q_1^{s-1} + \text{constant}$$

is homogeneous of degree  $p$  up to a constant, contradicting the hypothesis. Thus, if  $P$  has no homogeneous irreducible component of degree 1, we conclude that  $D^2P$  is negative definite on  $\mathcal{P}_n$ , so that  $P$  is strictly concave there; on the contrary, if it has some irreducible homogeneous component of degree 1, say,  $P_{I_\ell}$ , observe that if we fix all the other variables and if we evaluate  $P_{I_\ell}$  on the diagonal, we get an affine function. This achieves the proof.  $\square$

Now, we turn to optimization. For  $p > 0$ , let  $S_p^n = \{x \in \mathbf{R}^n; \sum_j |x_j|^p = 1\}$  denote the  $l^p$  unit sphere in  $\mathbf{R}^n$  and set  $S_{p,+}^n = \mathbf{R}_+^n \cap S_p^n$ . Our goal is to study the following problem:

*Given  $p > 0$  and a non-constant generalized polynomial with non-negative coefficients  $P$ , of degree  $h$  with  $h \leq p$ , characterize the argument(s) of*

$$\max_{x \in S_{p,+}^n} P(x). \quad (27)$$

It should be observed in the first place, since a generalized polynomial is continuous on  $\mathbf{R}_+^n$ , that the max in Problem (27) is indeed attained by compactness. We shall further investigate uniqueness and positivity properties of the argument of the max. In the second place, it is perhaps appropriate to

take a look at the limiting case  $p = \infty$  which was tacitly excluded in the statement of the problem. Then, it is clear that the maximum is attained by setting to 1 each variable which *actually appears* in  $P$  so that the problem is, in some sense, totally decoupled. For finite  $p$ , our first result concerns positivity:

**Proposition 2** *Assume a solution  $x^* = (x_1^* \cdots x_n^*)$  to Problem (27) satisfies  $x_i^* = 0$  for  $i \in I_1$  and  $x_i^* > 0$  for  $i \in I_2$ , where  $I_1$  and  $I_2$  partition the set of indices. Then one can decompose  $P$  as*

$$P(x) = P_1(x_{I_1}) + P_2(x_{I_2}), \quad (28)$$

where  $P_1$  is some (possibly zero) homogeneous generalized polynomial with non-negative coefficients of degree  $p$ .

*Proof:* Assume without loss of generality that  $I_1 = \{1 \cdots k\}$  and  $I_2 = \{k+1 \cdots n\}$ . Then  $k < n$  for  $x^* \neq 0$ . If a decomposition of the form (28) is impossible, this means that there exists a non-zero monomial in  $P$  whose degree with respect to the variables  $x_1, \dots, x_k$  is positive but less than  $p$ . Define  $\Phi_p = P \circ \varphi_p$ . This is a nonnegative generalized polynomial of degree  $h/p \leq 1$  deduced from  $P$  by dividing each exponent by  $p$ , say

$$\Phi_p(y) = \sum c_\beta y^\beta, \quad (29)$$

and it attains its maximum on  $S_{1,+}^n$  at

$$y^* = \varphi_p^{-1}(x^*) = (0, \dots, 0, y_{k+1}^*, \dots, y_n^*).$$

By a previous observation, one of the monomials, say,  $c_\gamma y^\gamma$  has non-zero degree less than 1 in the variables  $y_1 \cdots y_k$ . Thus, if we write

$$y^\gamma = y_1^{\gamma_1} \cdots y_n^{\gamma_n},$$

we have  $0 < \sum_{\ell=1}^k \gamma_\ell < 1$ . In particular, there is an index  $i \leq k$  such that  $0 < \gamma_i$  and

$$1 - \gamma_i > \sum_{1 \leq \ell \leq k, \ell \neq i} \gamma_\ell \stackrel{\text{def}}{=} \delta_i. \quad (30)$$

Select any  $j > k$  so that  $y_j^* > 0$  by definition of  $k$ . For  $0 \leq t \leq t_0 < y_j^*/k$ , the point

$$Y_t = (t, \dots, t, y_{k+1}^*, \dots, y_{j-1}^*, y_j^* - kt, y_{j+1}^*, \dots, y_n^*)$$

belongs to  $S_{1,+}^n$  and we can define  $G : [0, t_0] \rightarrow \mathbf{R}$  by the formula

$$G(t) = \Phi_p(Y_t).$$

This function is  $C^\infty$  for  $0 < t < t_0$ . Since every quantity involved is nonnegative, we have  $\partial \Phi_p / \partial \ell(Y_t) \geq 0$  for  $1 \leq \ell \leq n$ , and for  $\ell = i$  the stronger inequality:

$$\frac{\partial \Phi_p}{\partial i}(Y_t) \geq c_\gamma \gamma_i \frac{t^{\delta_i}}{t^{1-\gamma_i}} (y_{k+1}^*)^{\gamma_{k+1}} \dots (y_{j-1}^*)^{\gamma_{j-1}} (y_j^* - kt)^{\gamma_j} (y_{j+1}^*)^{\gamma_{j+1}} \dots (y_n^*)^{\gamma_n}. \quad (31)$$

Now, we evaluate

$$\frac{dG}{dt} = \sum_{\ell=1}^k \frac{\partial \Phi_p}{\partial \ell}(Y_t) - k \frac{\partial \Phi_p}{\partial j}(Y_t),$$

and we observe that the only negative contribution comes from the last term which is bounded for  $0 \leq t \leq t_0$  (since the only quantities appearing in the denominators of  $\partial \Phi_p / \partial j(Y_t)$  are of the form  $(y_j^* - kt)^{1-\beta_j}$ ), whereas the term corresponding to  $\ell = i$  is arbitrarily large when  $t$  is small enough by (31) and (30). Therefore, if  $t_0$  is small enough, we have  $dG/dt > 0$  hence (notice that the integral converges by continuity)

$$G(t_0) - G(0) = \int_0^{t_0} \frac{dG}{dt}(t) dt > 0,$$

contradicting the fact that  $y^*$  is a maximum.  $\square$

With the aid of Proposition 2, we can now treat uniqueness of the solution to (27) in an important special case.

**Theorem 3** *If  $P$  has no proper homogeneous irreducible component of degree  $p$  (in particular if  $P$  is irreducible), then Problem (27) has a unique solution  $x^*$ . Moreover,  $x^* \in \mathcal{P}_n$  in this case, and is the unique critical point of  $P$  on  $S_{p,+}^n \cap \mathcal{P}_n$ .*

*Proof:* by proposition 2, any solution  $x^*$  belongs to  $\mathcal{P}_n$ . Define again  $\Phi_p = P \circ \varphi_p$ , so that the maxima of  $\Phi_p$  on  $S_{1,+}^n$  are the images under  $\varphi_p^{-1}$  of those of  $P$  on  $S_{p,+}^n$ . Similarly, the critical points of  $\Phi_p$  on  $S_{1,+}^n \cap \mathcal{P}_n$  are the images under  $D\varphi_p^{-1}$  of those of  $P$  on  $S_{p,+}^n \cap \mathcal{P}_n$ . Since  $\Phi_p$  is a generalized polynomial of

degree at most 1, it is concave on  $\mathbf{R}_+^n$  by Theorem 2 and so is its restriction to the linear manifold  $S_{1,+}^n$ . Consequently its *maxima* on the latter form a convex set. Also, by concavity, these *maxima* coincide with the critical points of  $\Phi_p$  on  $S_{1,+}^n \cap \mathcal{P}_n$ , and the second derivative at such a point  $y^*$  is just the restriction of  $D^2\Phi_p(y^*)$  to the tangent space

$$\mathcal{T}_{y^*} S_{1,+}^n = \{x \in \mathbf{R}^n; \sum x_i = 0\}.$$

If  $P$  is reducible or if  $P$  is not homogeneous of degree  $p$ , the hypothesis strengthens to: “ $P$  has no irreducible homogeneous component of degree  $p$ ” so that  $\Phi_p$  has then no irreducible homogeneous component of degree 1; otherwise,  $P$  is irreducible homogeneous in degree  $p$  and so is  $\Phi_p$  in degree 1. According to each possibility, Theorem 2 tells us that the kernel of  $D^2\Phi_p(y^*)$  is either zero or one dimensional generated by  $y^*$ . In any case, this kernel intersects  $\mathcal{T}_{y^*} S_{1,+}^n$  at zero only. Therefore, the critical points of  $\Phi_p$  on  $S_{1,+}^n \cap \mathcal{P}_n$  are isolated, while at the same time forming a connected set since it is convex. Hence, there is a unique such point.  $\square$

To complete our study of Problem (27), we still have to examine what happens if  $P$  does have homogeneous irreducible components of degree  $p$ . To this effect, it will be convenient to generalize Problem (27) slightly and to consider

$$\max_{\substack{x \in \mathbf{R}_+^n \\ \|x\|_p = r}} P(x) \quad (32)$$

for  $r$  a non-negative real number. When  $r = 1$ , this is just Problem (27). Conversely, (32) reduces to (27) upon scaling each variable by  $r$ , so that the results established so far transpose immediately to Problem (32). In particular, if  $P$  has *no* proper homogeneous component of degree  $p$ , there is a unique argument for the max in (32) that we denote by  $x^*(r)$ . We have of course  $x^*(0) = 0$ . If  $r > 0$ , we know from Theorem 3 that  $x^*(r)$  is the unique critical point of  $P$  on  $\mathcal{P}_n \cap \{\|x\|_p = r\}$ ; this means that there exists a *Lagrange multiplier*  $\lambda^*(r)$  such that

$$\lambda^*(r) \frac{(x_i^*(r))^{p-1}}{r^{1-1/p}} = \frac{\partial P}{\partial i}(x^*(r)) \quad \forall i \in \{1, \dots, n\},$$

as this equation merely expresses that the gradient of  $P$  is proportional to the gradient of  $\|x\|_p$  at  $x^*(r)$ . Clearly,  $\lambda^*(r)$  is positive for  $x^*(r) \in \mathcal{P}_n$  and  $P$  is

not constant. Introducing the *Lagrangian* function

$$L_r(x, \lambda) = P(x) + \lambda(r - \|x\|_p),$$

this may be capsulized by saying that for  $r > 0$ , then  $(x^*(r), \lambda^*(r))$  is the unique critical point of  $L_r(x, \lambda)$  on  $\mathcal{P}_n \times \mathbf{R}$ . We shall need a few differential properties of  $x^*(r)$  and  $\lambda^*(r)$  as functions of  $r$ :

**Lemma 3** *If  $P$  has no proper homogeneous component of degree  $p$ , then, with the above notations,  $x^*(r)$  and  $\lambda^*(r)$  are  $C^\infty$  functions of  $r$  on  $(0, \infty)$ . If, in addition,  $P$  is not homogeneous of degree  $p$ , then*

$$\frac{d}{dr} [r^{1/p-1} \lambda^*(r)] \neq 0 \quad (33)$$

at every point of  $(0, \infty)$ .

*Proof:* put  $\Phi_p = P \circ \varphi_p$  and note that

$$(y^*(r), \mu^*(r)) = \left( \varphi_p^{-1}(x^*(r)), \frac{r^{1/p-1} \lambda^*(r)}{p} \right)$$

is the unique critical point over  $\mathcal{P}_n \times \mathbf{R}$  of the modified Lagrangian

$$L_r^1(y, \mu) = \Phi_p(y) + \mu(r^p - \sum_i y_i).$$

As in the proof of Theorem 3, we are now back to the elementary problem of maximizing a concave functional under some linear constraint, and the lemma is a standard application of the implicit function theorem granted Theorem 2 which guarantees nondegeneracy of the second derivative on the tangent space to the constraint. This computation we redo for the ease of the reader; by the implicit function theorem, we will obtain the desired smoothness if we show that the second derivative  $D^2 L_r^1$  is nonsingular at  $(y^*(r), \mu^*(r))$  because  $y^*$  and  $\mu^*$  will then be smooth functions of  $r$  and the same will obviously be true of  $x^*$  and  $\lambda^*$ . Compute this second derivative as

$$D^2 L_r^1(y^*(r), \mu^*(r)) = \begin{pmatrix} D^2 \Phi_p(y^*(r)) & -\mathbf{1} \\ -\mathbf{1}^T & 0 \end{pmatrix}, \quad (34)$$



where  $-\mathbf{1}$  stands for the vector in  $\mathbf{R}^n$  all components of which are -1's and where the superscript " $T$ " means "transpose". Assume  $(v, \nu) \in \mathbf{R}^n \times \mathbf{R}$  is in the kernel of this matrix. From the last row, we get  $\sum v_i = 0$  so that  $v$  belongs to the tangent space of  $S_{1,+}^n \cap \mathcal{P}_n$ . Then, multiplying (34) on the right by  $(v, \nu)$  and on the left by  $(v, \nu)^T$ , we deduce that  $v^T D^2 \Phi_p(y^*(r)) v = 0$ . But since  $\Phi_p$  has no proper irreducible homogeneous component of degree 1, we deduce from Theorem 2 (as in the proof of Theorem 3) that its second derivative restricted to the tangent space of  $S_{1,+}^n \cap \mathcal{P}_n$  is negative definite. Therefore, we have  $v = 0$  hence  $\nu = 0$  so that  $D^2 L_r^1(y^*(r), \mu^*(r))$  is non-singular. Now, equation (33) is equivalent to  $d\mu^*/dr \neq 0$ . Still from the implicit function theorem, and denoting by  $0$  the zero vector in  $\mathbf{R}^n$ , we have that

$$\frac{d\mu}{dr}(y^*(r), \mu^*(r)) = (\mathbf{0}^T \ 1) \left[ D^2 L_r^1(y^*(r), \mu^*(r)) \right]^{-1} \begin{pmatrix} \mathbf{0} \\ pr^{p-1} \end{pmatrix},$$

and this quantity cannot vanish if  $P$  has no homogeneous component of degree  $p$ , because it would mean that  $D^2 L_r^1(y^*(r), \mu^*(r))$  applied to some vector of the form  $(v, 0)^T$  yields  $(\mathbf{0}, pr^{p-1})^T$  so that  $D^2 \Phi_p(y^*(r)) v$  should in turn vanish and  $v$  itself should vanish since  $D^2 \Phi_p$  is non-singular by Theorem 3. This is absurd for  $r > 0$ .  $\square$

For our purposes, we need to know a bit more about the behaviour of the objective function in Problem (32), that we define as

$$M_{P,p}(r) = \max_{\substack{x \in \mathbf{R}_+^n \\ \|x\|_p = r}} P(x). \quad (35)$$

For instance, when  $P$  is homogeneous of degree  $h$ , then  $M_{P,p}$  is just  $r^h M_{P,p}(1)$ . When  $P$  is not homogeneous, things get more complicated but the properties of  $M_{P,p}$  that we will use are gathered in the following lemma.

**Lemma 4** *If  $P$  has no homogeneous component of degree  $p$ , then  $M_{P,p}$  is a continuous function on  $\mathbf{R}_+$  which is  $C^\infty$  on  $(0, \infty)$ . Moreover,  $M_{P,p}(t^{1/p})$  is a strictly concave function on  $(0, \infty)$  whose first derivative is positive and whose second derivative is negative there. When  $p \leq 1$ , the same is therefore true of the function  $M_{P,p}(r)$  itself.*

*Proof:* smoothness of  $M_{P,p}$  follows from Lemma 3 and from the relation

$$M_{P,p} = L_r(x^*(r), \lambda^*(r)). \quad (36)$$

Continuity of  $M_{P,p}$  at  $0^+$  is obvious. Note also that

$$\frac{d M_{P,p}}{dr}(r) = \lambda^*(r), \quad (37)$$

which is an ultraclassical result in optimization asserting that the Lagrange multiplier can be interpreted as the sensitivity of the *optimal* value to the constraint level; (37) drops out immediately from (36) and from the fact that  $(x^*(r), \lambda^*(r))$  is critical for  $L_r$ . From (37), we see that  $dM_{P,p}/dr > 0$  on  $(0, \infty)$ , hence also  $dM_{P,p}(t^{1/p})/dt > 0$  by the chain rule. Setting as usual  $\Phi_p = P \circ \varphi_p$ , we readily observe that

$$M_{P,p}(t^{1/p}) = M_{\Phi_p,1}(t), \quad (38)$$

and it follows then from (37) and (33) that  $d^2 M_{\Phi_p,1}/dr^2$  is never zero. It will in fact be negative for the function is concave; indeed, let  $x^*$  and  $y^*$  have  $l^1$  norm  $r_1$  and  $r_2$  respectively, and be such that  $\Phi_p(x^*) = M_{\Phi_p,1}(r_1)$  and  $\Phi_p(y^*) = M_{\Phi_p,1}(r_2)$ . Since  $\Phi_p$  has degree at most 1, it is concave on  $\mathbf{R}_+^n$  by Theorem 2. Hence, we get for  $\mu_1 \geq 0$  and  $\mu_2 \geq 0$  satisfying  $\mu_1 + \mu_2 = 1$ :

$$\mu_1 M_{\Phi_p,1}(r_1) + \mu_2 M_{\Phi_p,1}(r_2) = \mu_1 \Phi_p(x^*) + \mu_2 \Phi_p(y^*) \leq \Phi_p(\mu_1 x^* + \mu_2 y^*).$$

As  $\|\mu_1 x^* + \mu_2 y^*\|_1 \leq \mu_1 r_1 + \mu_2 r_2$ , it follows that

$$\Phi_p(\mu_1 x^* + \mu_2 y^*) \leq M_{\Phi_p,1}(\mu_1 r_1 + \mu_2 r_2),$$

for  $M_{\Phi_p,1}$  is an increasing function. Finally, when  $p \leq 1$ , it is immediate from the chain rule that  $dM_{P,p}(t^{1/p})/dt > 0$  and  $d^2 M_{P,p}(t^{1/p})/dt^2 < 0$  together imply  $d^2 M_{P,p}(t)/dt^2 < 0$ . This achieves the proof.  $\square$

We can now assess positivity and uniqueness in the general case for Problem (27). *For definiteness, we shall assume that  $P$  has no constant coefficient.* This does not change the arguments of the *maximum* in Problem (27) and consequently does not affect our results. But it is to the effect that the irreducible components of  $P$  will in turn be free of constant coefficients, hence homogeneous components will really be homogeneous, and not only up to an additive

constant. So, assume that

$$P(x) = \sum_{\ell=1}^k P_{I_\ell}(x_{I_\ell}) \quad (39)$$

is the irreducible decomposition of  $P$ , each  $I_\ell$  being of cardinality  $n_\ell$ , and, say  $P_{I_1}, \dots, P_{I_s}$  are homogeneous of degree  $p$ . Define  $H = \cup_{\ell=1}^s I_\ell$ ,  $K = \cup_{\ell=s+1}^k I_\ell$ . Reordering the indices is necessary, we may suppose that  $\{1, \dots, m\}$  are those indices  $\ell \leq s$  for which  $M_{P_{I_\ell}, p}(1)$  is *maximum*, and we denote this *maximum* common value by  $M_h$ . This gives rise to a partition  $H = H_1 \cup H_2$  of  $H$ , with

$$H_1 = \cup_{\ell=1}^m I_\ell \quad \text{and} \quad H_2 = \cup_{\ell=m+1}^s I_\ell.$$

We set accordingly

$$P_{H_1}(x_{H_1}) = \sum_{\ell=1}^m P_{I_\ell}(x_{I_\ell}), \quad P_{H_2}(x_{H_2}) = \sum_{\ell=m+1}^s P_{I_\ell}(x_{I_\ell}),$$

$$\text{and } P_K(x_K) = \sum_{\ell=s+1}^k P_{I_\ell}(x_{I_\ell}),$$

which depend on

$$n_{H_1} = \sum_{\ell=1}^m n_\ell, \quad n_{H_2} = \sum_{\ell=m+1}^s n_\ell, \quad \text{and } n_K = \sum_{\ell=s+1}^k n_\ell$$

variables respectively. These are generalized polynomial with non-negative coefficients, the first two being homogeneous of degree  $p$  while the third is of degree at most  $p$  and has *no* homogeneous irreducible component of degree  $p$ . For  $1 \leq \ell \leq m$ , we shall denote by  $z_{I_\ell}^* \in \mathcal{P}_{n_\ell} \cap S_{p,+}^{n_\ell}$  the maximizing vector such that  $P_{I_\ell}(z_{I_\ell}^*) = M_h$ , whose existence and uniqueness is asserted in Theorem 3. If  $s < k$ , that is, if  $P$  is not homogeneous of degree  $p$ , we further define  $z_K^*(r) \in \mathcal{P}_{n_K} \cap \{\|x\|_p = r\}$  to be the solution to Problem (32) for  $P_K$ , whose existence and uniqueness again follows from Theorem 3, and by

$$\lambda_K^*(r) = r^{1-1/p} (z_K^*)_i^{1-p}(r) \frac{\partial P_K}{\partial i}(z_K^*(r))$$

the associated Lagrange multiplier as introduced in Lemma 3 (whose value does not depend on  $i \in \{1, \dots, n_K\}$ ). We simply set  $z_K^* = z_K^*(1)$  and  $\lambda_K^* = \lambda^*(1)$  for the pair associated with Problem (27). When  $s = k$ , then  $K$  is of course empty so we should forget about  $P_K$ ,  $x_K^*$ ,  $z_K^*$ , and  $\lambda_K^*$  in the statement of the next theorem.

**Theorem 4** *With the above notations, the set of solutions to Problem 27 consists of those  $x^*$  satisfying*

$$\begin{aligned} x_{I_\ell}^* &= \mu_\ell z_{I_\ell}^* \quad \text{for } 1 \leq \ell \leq m \quad \text{and } \mu_\ell \in \mathbf{R}^+ \quad \text{subject to } \sum_{\ell=1}^m \mu_\ell^p = 1 - r_0^p, \\ x_{H_2}^* &= 0, \quad x_K^* = z_K^*(r_0), \end{aligned} \tag{40}$$

where  $r_0$  is the unique maximum on  $[0, 1]$  of the function

$$g(r) = (1 - r^p)M_h + M_{P_K, p}(r).$$

We have  $0 < r_0 < 1$  unless  $s < k$  and  $\lambda_K^* \geq pM_h$ , in which case  $r_0 = 1$  and  $x_{H_1}^* = 0$ . In particular, the solution  $x^*$  is unique if, and only if, either  $m = 1$  or  $s < k$  and  $\lambda_K^* \geq pM_h$  (these two cases are not exclusive). There exists a solution in  $\mathcal{P}_n$  if and only if  $m = s$  (i. e.  $H_2$  is void). Every solution lies in  $\mathcal{P}_n$  if, and only if,  $m = s = 1$  (hence the solution is unique) and  $\lambda^* < pM_h$  in case  $s < k$ ; then, the solution  $x^*$  is the unique critical point of  $P$  on  $S_{p,+}^n \cap \mathcal{P}_n$ .

*Proof:* let  $x^*$  be a solution; then the optimal value is equal to

$$M_{P,p}(1) = M_{P_{H_1},p}(\|x_{H_1}^*\|_p) + M_{P_{H_2},p}(\|x_{H_2}^*\|_p) + M_{P_K,p}(\|x_K^*\|_p),$$

so, by homogeneity of  $P_{H_1}$  and  $P_{H_2}$ ,

$$M_{P,p}(1) = M_h \|x_{H_1}^*\|_p^p + \sum_{\ell=m+1}^s M_{P_{I_\ell},p} \|x_{I_\ell}^*\|_p^p + M_{P_K,p}(\|x_K^*\|_p). \tag{41}$$

As  $M_h > M_{P_{I_\ell},p}$  for  $m+1 \leq \ell \leq s$ , and since we are maximizing under the constraint

$$\sum_{\ell=1}^k \|x_{I_\ell}\|_p^p = 1,$$

it is clear that  $x_{H_2}^*$  must be zero and that  $x_{I_\ell}^* = \mu_\ell z_{I_\ell}^*$  for  $1 \leq \ell \leq m$ , where the  $\mu_\ell$ 's should satisfy  $\sum_{\ell=1}^m \mu_\ell^p = 1 - \|x_K^*\|_p^p$  but otherwise produce the same value for  $M_{P,p}(1)$ . If  $s = k$ , then  $K$  is void,  $M_{P_K,p}$  is not present in (41), and  $\|x_K^*\|_p$  should be interpreted as zero. Since also  $r_0 = 0$  in this case, (40) then holds true. If  $s < k$ , (41) reads now:

$$M_{P,p}(1) = M_h(1 - \|x_K^*\|_p^p) + M_{P_K,p}(\|x_K^*\|_p) = g(\|x_K^*\|_p).$$

Setting  $t = \|x_K^*\|_p^p$ , we get

$$g(\|x_K^*\|_p) = M_h(1 - t) + M_{P_K,p}(t^{1/p}),$$

and, since  $P_K$  has no homogeneous components of degree  $h$ , we deduce from Lemma 4 that the above is a smooth strictly concave function of  $t$  on  $[0, \infty)$ . Therefore, it has a unique *maximum* on  $[0, 1]$ , which is necessarily attained at  $t_0 = \|x_K^*\|_p^p$  by the optimality of  $x^*$ . From Proposition 2, we also see that  $t_0 \neq 0$  for otherwise  $P$  would reach a *maximum* on  $S_{p,+}^n$  at the point  $x^*$  satisfying  $x_K^* = 0$ , whereas  $P_K$  is a sum of components, none of which is homogeneous of degree  $p$ . This shows in particular that  $x_K^* = z_K^*(r_0)$ , where  $r_0 = t_0^{1/p}$  is indeed the *maximum* of  $g(r)$  on  $(0, 1)$ . Now, this *maximum* is attained on  $(0, 1]$ , and  $x_{H_1}$  will be zero if, and only if, it is attained at 1, that is, if, and only if,

$$\frac{g(t^{1/p})}{dt}(1) \geq 0.$$

Expanding the derivative using (37) yields then  $\lambda_K^* \geq pM_h$ . The remaining assertions are now obvious.  $\square$

### Remarks

1) In principle, Theorem 4 reduces the general case of Problem (27) to a sequence of situations covered by Theorem 3 where the *optimum* can be computed by almost any method in optimization since, composing with  $\varphi_p$ , we are back to maximizing a smooth strictly concave function over a convex open subset of a linear space and we know the *maximum* is attained. In particular, the problem of deciding which variables are zero at an *optimum* is equivalent to determining the homogeneous irreducible components of degree  $p$ —a combinatorial step— and then comparing the optimal values they achieve on

Problem (27) –an analytical step– while these values can be computed rather easily as we just mentioned.

2) As we noticed already, a point  $x^c \in S_{p,+}^n \cap \mathcal{P}_n$  is critical for  $P$  if and only if there exists a Lagrange multiplier  $\lambda$  such that

$$\lambda (x_i^c)^{p-1} = \frac{\partial P}{\partial i}(x^c) \quad \forall i \in \{1, \dots, n\}. \quad (42)$$

Now, as in the proof of Theorem 2, let  $\alpha_m$  be the smallest non-zero exponent for a variable appearing in  $P$ ; changing  $P$  into  $P \circ \varphi_{\alpha_m}$  if necessary, we may assume in Problem (27) that  $\alpha_m \geq 1$  implying  $p \geq 1$  also. In this case, equation (42) makes sense for any point in  $S_{p,+}^n$ , that is even if some components of  $x^c$  are equal to zero, so that we could *define* a critical point of  $P$  on  $S_{p,+}^n$  in this way. If  $\alpha_m = p = 1$ , then  $P$  is a linear polynomial and the *maxima* of  $P$  are generally not critical points. But if  $p > 1$ , they are critical because if  $x^*$  is such a point and  $I_1, I_2$  partition the set of indices in such a way that  $x_{I_1}^* = 0$  while  $x_i^* > 0$  for  $i \in I_2$ , then we know from Proposition 3 that  $P$  decomposes as  $P_1(x_{I_1}) + P_2(x_{I_2})$  where  $P_1$  is homogeneous of degree  $p$  and those equations in (42) that correspond to null components of  $x^*$  can be read  $0 = 0$  so that they are automatically satisfied. It is natural to ask for the converse, namely is a critical point necessarily a *maximum*? The answer is no: for  $p \geq 3$ , if we denote by  $(x_1^c, x_2^c)$  the *maximum* of  $x_1 x_2$  on  $S_{p,+}^2$ , the point  $(x_1^c, x_2^c, 0, 0)$  is critical for  $P(x) = x_1 x_2 + x_2 x_3 x_4$  on  $S_{p,+}^4$  but is not a *maximum* for it does not belong to  $\mathcal{P}_4$  though  $P$  is irreducible.

*To recap, assuming  $\alpha_m \geq 1$ , we have in Problem (27) that a critical point lying in  $\mathcal{P}_n$  is necessarily a solution whereas a critical point with some zero components may not be a solution.*

As a byproduct of the preceding discussion, we may also point out a kind of non-linear generalization of the Perron-Frobenius theorem [6] in case  $A$  is symmetric.

**Corollary 1** *Let  $A$  be a real symmetric  $n \times n$  matrix with non-negative entries. For any  $\alpha \geq 1$ , there exists a nonzero  $x^* \in \mathbf{R}_+^n$  such that*

$$Ax^* = \lambda^* \varphi_{1/\alpha}(x^*), \quad \text{for some } \lambda^* \geq 0. \quad (43)$$

*If  $A$  is irreducible, then  $x^*$  is unique up to a multiplicative constant and belongs to  $\mathcal{P}_n$ . If  $A$  is reducible, the solution is no longer unique. If  $\alpha$  is an integer (so*

that equation (43) makes sense for negative  $x$ 's as well), the largest possible value for  $\lambda^*$  is also the largest value of  $|\lambda|$  for which

$$Ax = \lambda \varphi_{1/\alpha}(x) \quad (44)$$

is solvable with respect to  $x \in \mathbf{R}^n - \{0\}$  for some (possibly negative)  $\lambda$ . When  $\alpha > 1$  (the non-linear case), the vector  $x^*$  associated with this largest value belongs to  $\mathcal{P}_n$ .

*Proof:* if  $A = 0$ , there is nothing to prove. Otherwise, set  $P(x) = x^T Ax$ , which is homogeneous of degree 2. Consider Problem (27) with  $p = \alpha + 1$ . Let  $I_1 \cdots I_k$  denote the irreducible partition of  $A$ , so that equation (43) splits into  $k$  subequations in the  $x_{I_j}$ 's (because  $A$  is symmetric). Now, the irreducible decomposition of  $P$  writes

$$P(x) = \sum_{\ell=1}^k P_{I_\ell}(x_{I_\ell}), \quad (45)$$

each  $P_{I_\ell}$  being again homogeneous of degree 2. Let  $x^*$  be a solution to (27). By proposition 2, the indices of the null coordinates of  $x^*$  range over a union  $\cup_{j \in J} I_j$ , for some proper subset  $J$  of  $\{1 \cdots k\}$ , and this union can be nonempty only when  $p = 2$ , that is, when  $\alpha = 1$ . In this case, the subequations of (43) corresponding to the  $x_{I_j}$ 's for  $j \in J$  read  $0 = 0$  and are automatically satisfied; the problem then reduces to a similar one in fewer variables. Altogether, we may assume  $x^* \in \mathcal{P}_n$ . Then, writing that  $x^*$  is a critical point of  $P$  on  $S_{p,+}^n \cap \mathcal{P}_n$  and observing that the vector  $\varphi_{1/(p-1)}(x^*)$  is normal to  $S_{p,+}^n$  at  $x^*$ , we get  $Ax^* = \lambda^* \varphi_{1/(p-1)}(x^*)$  for some  $\lambda^*$  which is obviously nonnegative, that is, (43) is satisfied. If  $A$  is irreducible, every non-zero solution to (43) belongs to  $\mathcal{P}_n$  (easy checking) and Theorem 3 tells us that there is a unique critical point. This means that a solution of unit  $l^p$  norm to (43) is unique in this case. If  $A$  is reducible, a solution is no longer unique as it is clear from what precedes that we may set  $x_{I_\ell}$  to zero for  $\ell$  ranging over a strict subset  $\{1, \dots, k\}$ , and still find a nonzero solution in terms of the remaining variables (of course this will not, in general, lead to a solution of (43) which is at the same time a *maximum* of  $P$  on  $S_{p,+}^n$ ). Finally, if  $\alpha$  is an integer and  $x \in \mathbf{R}^n$  any non-zero solution to (44) of unit  $l^p$  norm, we have  $x^T Ax = \lambda$  and, since  $A$  is non-negative, this number cannot exceed the *maximum* of  $P$  on  $S_{p,+}^n$ .  $\square$

### Remarks

- i) Extending a remark of [8] concerning this theorem, we may notice that *existence* of a solution to (43) would also follow from Brouwer's fixed-point theorem as applied to the map  $x \rightarrow \varphi_\alpha(Ax)/\|\varphi_\alpha(Ax)\|_p$  from  $S_{2,+}^n$  into itself.
- ii) Corollary 1 would remain valid for matrices depending on  $x$ , provided  $Ax$  is the gradient of some non-negative generalized polynomial. This entails algebraic conditions that we shall not analyse here.

## 6 A generalization to product-type constraints

Some of the previous results extend easily to the case where the set of indices  $I = \{1, \dots, n\}$  is partitioned into  $d$  blocks  $J_1, \dots, J_d$  of respective sizes  $\nu_1, \dots, \nu_d$ , that is

$$J_i = \{x_k; \sum_{\ell=1}^{i-1} \nu_\ell < k \leq \sum_{\ell=1}^i \nu_\ell\}$$

(the empty sum occurring if  $i = 1$  has of course to be interpreted as zero), with  $\sum_i \nu_i = n$  and the constraint is now  $\|x_{J_i}\|_p = 1$  for each  $i \in \{1, \dots, d\}$ . We shall not analyse the solution to this generalized problem as completely as we did for Problem (27). However, since this generalization is relevant to the applications presented in the introduction, we shall proceed in this section with those results that can be stated in a fairly general manner and at the same time warrant the search for critical points in practice. Letting  $\nu$  stand for the  $d$ -tuple  $(\nu_1, \dots, \nu_d)$ , the set over which we optimize becomes

$$\mathcal{S}_{p,+}^\nu \stackrel{\text{def}}{=} S_{p,+}^{\nu_1} \times S_{p,+}^{\nu_2} \times \dots \times S_{p,+}^{\nu_d},$$

and we state the problem formally as:

*Given  $p > 0$  and a non-constant generalized polynomial  $P$  with non-negative coefficients of degree  $h$  with  $h \leq p$ , characterize the argument(s) of*

$$\max_{x \in \mathcal{S}_{p,+}^\nu} P(x). \tag{46}$$

Proposition 2 carries over *mutatis mutandis* to Problem (46):



**Proposition 3** *Assume a solution  $x^*$  to Problem (46) satisfies  $x_i^* = 0$  for  $i \in I_1$  and  $x_i^* > 0$  for  $i \in I_2$ , where  $I_1$  and  $I_2$  partition the set of indices. Then one can decompose  $P$  as*

$$P(x) = P_1(x_{I_1}) + P_2(x_{I_2}), \quad (47)$$

where  $P_1$  is some (possibly zero) homogeneous generalized polynomial with non-negative coefficients of degree  $p$ .

*Proof:* the proof is similar to that of proposition 2 except for two facts:

1) we cannot assume that  $I_1 = \{1, \dots, k\}$ , because this time the indices have been fixed by the way we formulated the constraints. This creates only notational inconvenience.

2) the variable  $y_j$  such that  $y_j^* > 0$  has to be replaced by a collection  $y_\alpha$  with  $\alpha \in \Lambda$ , where  $y_\alpha^* > 0$  and  $\Lambda$  contains exactly one element in each non-empty intersection  $J_i \cap I_1$  as  $i$  ranges over  $\{1, \dots, d\}$ . For  $\alpha \in J_i \cap I_1$ , we then define  $n_\alpha > 0$  to be the cardinality of  $J_i \cap I_1$ .

Letting now  $Y_t$  be the vector such that each component of  $(Y_t)_{I_1}$  is  $t$  and  $(Y_t)_\alpha = y_\alpha^* - n_\alpha t$  while all other components of  $Y_t$  are equal to those of  $y^*$ , we leave it to the reader to check that the proof of Proposition 2 carries over with obvious changes.  $\square$

We now obtain a straightforward generalization of Theorem 3:

**Theorem 5** *The solutions to Problem (46) form a nonempty, connected, and closed subset of  $\mathcal{S}_{p,+}^\nu$ . The set of critical points and the set of maxima of  $P$  coincide on  $\mathcal{S}_{p,+}^\nu \cap \mathcal{P}_n$ . If  $P$  has no proper irreducible homogeneous components of degree  $p$  (in particular if  $P$  is irreducible), then there is a unique solution to Problem (46) and it lies on  $\mathcal{S}_{p,+}^\nu \cap \mathcal{P}_n$  where it is thus the unique critical point of  $P$ . More generally, if  $x^*$  denotes any solution to Problem (46), the components of  $x^*$  that are not involved in a proper irreducible homogeneous component of degree  $p$  are uniquely determined.*

*Proof:* we mimic the proof of Theorem 3. The set of solutions is nonempty and closed by the compactness of  $\mathcal{S}_{p,+}^\nu$  and the continuity of  $P$ . Put  $\Phi_p = P \circ \varphi_p$ . Since  $\Phi_p$  has degree at most 1, it is concave on  $\mathbf{R}_+^n$  by Theorem 2 and so is its restriction to the convex subset  $\mathcal{S}_{1,+}^\nu$  of the affine subspace

$$\{y \in \mathbf{R}^n; \sum_{j \in J_i} y_j = 1 \text{ for } 1 \leq i \leq d\};$$

consequently, the *maxima* of  $\Phi_p$  on  $\mathcal{S}_{1,+}^\nu$  form a convex set which is therefore connected and since they are the images under  $\varphi_p^{-1}$  of the *maxima* of  $P$  on  $\mathcal{S}_{p,+}^\nu$ , the latter must form a connected set as well. Further, we get by concavity that any critical point of  $\Phi_p$  on the linear manifold  $\mathcal{S}_{1,+}^\nu \cap \mathcal{P}_n$  is a *maximum* of  $\Phi_p$  with respect to  $\mathcal{S}_{1,+}^\nu$ , and conversely any *maximum* lying in  $\mathcal{P}_n$  is a critical point. Since the critical points of  $\Phi_p$  on  $\mathcal{S}_{1,+}^\nu \cap \mathcal{P}_n$  are the images under  $D\varphi_p^{-1}$  of those of  $P$  on  $\mathcal{S}_{p,+}^\nu \cap \mathcal{P}_n$ , we see that critical points and *maxima* of  $P$  coincide on  $\mathcal{S}_{p,+}^\nu \cap \mathcal{P}_n$ .

If  $P$  has no proper irreducible component which is homogeneous of degree  $p$ , Proposition 3 tells us that the solutions to Problem (46) belong to  $\mathcal{P}_n$  and the same then holds for any *maximum*, say  $y^*$ , of  $\Phi_p$  on  $\mathcal{S}_{1,+}^\nu$ . We get, as in the proof of Theorem 3, that the kernel of  $D^2\Phi_p(y^*)$  is either zero or one dimensional generated by  $y^*$ , and therefore intersects

$$\mathcal{T}_{y^*} \mathcal{S}_{1,+}^\nu = \{y \in \mathbf{R}^n; \sum_{j \in J_i} y_j = 0 \text{ for } 1 \leq i \leq d\} \quad (48)$$

at zero only. Therefore, the critical points of  $\Phi_p$  on  $\mathcal{S}_{1,+}^\nu \cap \mathcal{P}_n$  are isolated, and the same is true of those of  $P$  on  $\mathcal{S}_{p,+}^\nu \cap \mathcal{P}_n$ . Because we just proved they form a connected set, there is only one such point. More generally, if we assume that  $\Phi_p$  assumes a *maximum* both at  $y^*$  and  $z^*$ , then  $\Phi_p(ty^* + (1-t)z^*)$  is also *maximum* for  $0 \leq t \leq 1$  so the vector  $y^* - z^*$  lies in the kernel of  $D^2\Phi_p(y^*)$ , and this implies by theorem 2 that the coordinates of  $y^* - z^*$  whose index is not involved in some homogeneous component of degree 1 do vanish.  $\square$

**Remark** The second remark we made after Theorem 4 remains of course valid: if the smallest exponent for a variable in  $P$  is not less than 1, we can define critical points on  $\mathcal{S}_{p,+}^\nu$ . However, such a point need not be a *maximum* if it does not belong to  $\mathcal{P}_n$ .

At this point, it would be possible to derive an analog of Theorem 4 for Problem (46) by comparing the optimal values of the homogeneous and the remaining parts of  $P$ . However, it seems hardly worthwhile to build such a general statement because it would be rather intricate. Indeed, the function  $M_{P,p}(r)$  defined in (35) should be replaced by some  $M_{P,p}(r_1, \dots, r_d)$  whose behaviour is more complex even if  $P$  is homogeneous because it need not be homogeneous with respect to each  $x_{J_i}$  separately. Leaving it to the interested reader to analyse further specific cases, we shall rather illustrate how multiple

constraints may interact with the irreducible decomposition by giving a simple triangular criterion for uniqueness.

**Proposition 4** *Let  $P$  have  $s \geq 1$  homogeneous irreducible components of degree  $p$  and let  $I_1, \dots, I_s$  be the corresponding elements of the irreducible partition. If the ordering  $1, \dots, s$  can be arranged so that for each  $I_j$  there is a  $i(j) \in \{1, \dots, d\}$  with the property that  $J_{i(j)} \cap I_j \neq \emptyset$  but  $J_{i(j)} \cap I_k = \emptyset$  for  $k < j$ , then the solution to Problem (46) is unique.*

**Remark** When  $P$  is irreducible, this is nothing new in view of Theorem 5. We can even assert in this case that the solution is the unique critical point on  $\mathcal{S}_{p,+}^\nu \cap \mathcal{P}_n$ .

*Proof:* let  $x^*$  be a solution and assume first that  $x^* \in \mathcal{P}_n$ . Define as usual  $\Phi_p = P \circ \varphi_p$ , so that  $y^* = \varphi_p^{-1}(x^*)$  is a critical point of  $\Phi_p$  on  $\mathcal{S}_{1,+}^\nu$ . It follows from Theorem refconcave1 that the kernel of  $D^2\Phi_p(y^*)$  consists of vectors of the form  $\sum_j \mu_j y_{I_j}^*$  with  $\mu_j \in \mathbf{R}$ . By (48), such a vector belongs to the tangent space of  $\mathcal{S}_{1,+}^\nu$  if and only if

$$\sum_{j=1}^s \mu_j \sum_{\ell \in J_i \cap I_j} y_\ell = 0 \text{ for } 1 \leq i \leq d\}.$$

This means that the vector  $\mu = (\mu_1, \dots, \mu_s)^T$  lies in the kernel of the  $d \times s$  matrix whose  $(i, j)$  entry is  $\sum_{\ell \in J_i \cap I_j} y_\ell$ . By our hypothesis, this matrix has a nonsingular triangular submatrix of size  $s$  so that  $\mu = 0$  and  $D^2\Phi_p(y^*)$  restricted to  $\mathcal{T}_{y^*} \mathcal{S}_{1,+}^\nu$  is nonsingular. Consequently,  $y^*$  is isolated and so is  $x^*$ . The latter is therefore unique by Theorem 5.

Assume now that  $x^*$ , hence also  $y^*$ , has some zero components whose indices then range over a union  $\cup I_\alpha$  by Proposition 3. If  $z^* \neq y^*$  is another point at which  $\Phi_p$  attains a *maximum* on  $\mathcal{S}_{1,+}^\nu$ , the zero components of  $z^*$  also have indices ranging over a union  $\cup I_\beta$ . By concavity,  $\Phi_p$  attains a *maximum* at every point of the form  $ty^* + (1-t)z^*$  with  $0 < t < 1$ , and the zero components of these points have indices ranging over  $\cup I_\alpha \cap \cup I_\beta$  which is again a union, say,  $\cup I_\gamma$ . Setting  $y_{I_\gamma} = 0$  for each  $\gamma$ , we get from  $\Phi_p(y)$  a new generalized polynomial in fewer variables which still meets the assumptions (because we have merely suppressed a few homogeneous components) but attains a *maximum* at infinitely many points with positive components. This contradicts the first part of the proof and establishes the Proposition.  $\square$

## 7 Finding the maximal mode of a Potts model

In this section, we study in more details the behaviour of DPA on the optimization of randomly weighted *square* Markov Random Fields. The idea is to generate reasonably large graphs, so that it is possible to make an exhaustive search of the optimal solution, and so perform a more objective evaluation of the results of DPA.

We build such an MRF in the following way. The sites are the pixels on an image, and the cliques are determined by the maximum distance of neighbouring pixels. This way, we successively study cliques of order 2, generated by 4-connectivity, then cliques of order 4, generated by 8-connectivity. We also vary the number of labels, or states, for each site, from 2 (corresponding to the Ising model), to 4 (an instance of a Potts model). Once the type of connectivity, the number of labels and the size of the graph (i.e. the image) have been selected, then the values of the clique potentials are generated randomly (typically with uniform distribution between 0 and 1, but the results do not change much if another distribution is used). In each case, 10 and 20 such random MRF's are generated. The value of the optimal configuration can be determined exactly by dynamic programming, following Derin & Elliott [4]. If the image has width  $w$  and height  $h$ , then the complexity of the search is  $hw^M$ . Images of width up to 10, and height up to 20 can thus easily be searched on a workstation. The experiments related here were made on  $5 * 5$  to  $16 * 32$  images.

Applying DPA to this problem is straightforward. The maximization of  $f$  on  $\mathcal{K}_p$  can be performed by any method as long as  $p \geq \deg(G)$ , for we know that the maximum is unique in this case. As suggested in [1], a direct generalization of the iterative power method for finding the eigenvector corresponding to the maximum eigenvalue can be applied. It consists of selecting some  $X^{(0)}$  and subsequently compute

$$\forall i : X_i^{(n+1)} = \alpha_i^{(n+1)} \left( D_{X_i} f(X_i^{(n)}) \right)^{\frac{1}{p-1}}, \quad (49)$$

where  $X_i$  denotes the vector  $(x_{i,1} \cdots x_{i,M})$  and  $D_{X_i} f$  the partial gradient with respect to these variables, while the  $\alpha_i^{(n+1)}$ 's are adjusted so that each  $X_i$  has unit  $l^p$ -norm.

This simply means that, at each iteration, we select on the pseudo-sphere of degree  $p$  the point where the normal is parallel to the gradient of  $f$ . Obviously, the unique fixed point is the *maximum* we are looking for. It has not yet been established that the sequence of values generated in this way is monotonically increasing although experiments indicate that this is always the case in practice. Besides, convergence is fast. It is anyway easy to check at each iteration whether  $f$  increases, and if not to fall back on a standard gradient. This was never necessary in all the applications we tried. The reason for that is still an open question.

Another important remark is that, for some unknown  $p_b < \deg(C)$ , a bifurcation occurs and the maximum is no longer unique. Here comes the heuristic part of the algorithm: we simply ignore the bifurcation, and carry on with the iterative power algorithm even  $p < p_b$ . More precisely, the complete algorithm goes as follows:

1. set  $X^{(0)} = [1, \dots, 1], k = 1$ ,
2. while  $p > 1$  do
  - $Y^{(k)} \leftarrow$  the projection of  $X^{(k-1)}$  on  $\mathcal{K}_p$ ,
  - find  $X^{(k)} = \max_{\mathcal{K}_p} f(X)$  using the iterative power method, starting from  $Y^{(k)}$ ,
  - $k \leftarrow k + 1, p \leftarrow p - \beta$
  - od

Here, projecting a non-negative  $Y$  on  $\mathcal{K}_p$  is performed by scaling separately for each site, i.e.  $X_i = \alpha_i Y_i, \forall i$ , just as shown above.

For  $p < p_b$ , it could, theoretically, happen that  $Y^{(k)}$  is a critical point, possibly a local minimum, or a saddle-point, in which case the iterative power method gets stuck. Needless to say this is very unlikely, and never occurred in practice.

We have also run comparisons with a Gibbs sampler, setting the number of iterations for it to the total number of iterations (summed over the  $\beta$ 's) set for DPA. The results are displayed on Tables 3 (for 4-connectivity and cliques

$N_\beta$	$E_{DPA}$	$V_{DPA}$	$V_{opt}$
2	2.2	128.40	128.74
3	2.6	128.43	128.74
5	2.6	128.45	128.74

Table 1: Size  $8 \times 8$ , 2 labels, 4-connectivity,  $N_{it} = 50$ ,  $Th = 10^{-5}$ 

$Th$	$E_{DPA}$	$V_{DPA}$	$V_{opt}$
$10^{-7}$	3.6	128.25	128.74
0.3	3.6	128.19	128.74

Table 2: Size  $8 \times 8$ , 2 labels, 4-connectivity,  $N_\beta = 3$ ,  $N_{it} = 5$ 

of order 2) and 4 (for 8-connectivity and cliques of order 4). Here,  $E_G$  is the average number of errors for Gibbs, and  $V_G$  the average value of the criterion. The last column (*DPA/Gibbs*) displays the percentage of trials for which the value reached by DPA was better than the value reached by Gibbs.

The results on Table 3 show that, for 2 labels and cliques of order 2 (4-connectivity), DPA is definitely better, and quite close to optimal. For 3 or more labels, DPA is still better than Gibbs, but farther from optimal. For higher-order cliques (Table 4), DPA degrades faster than Gibbs.

We have also applied on larger graphs a search by dynamic programming with pruning of the current hypotheses, thus implementing a variation of the Viterbi algorithm. On an image 8-pixel wide, and 2 labels, an exhaustive search implies to maintain  $2^9$  current hypotheses. We have found, as shown on Table 5, that the results are quite poor as soon as pruning (by discarding the worst hypothesis) exceeds 50%. This makes this last method just as intractable as optimal search.

## 8 Conclusion

Having analysed in detail the  $l_p$ -constrained maximization problem for generalized polynomials of degree at most  $p$ , it is natural to ask about the best

$N_{lab}$	$E_{DPA}$	$E_G$	$V_{DPA}$	$V_G$	$V_{opt}$	$DPA/Gibbs$
2	12.1	29.4	211	206	214.2	100%
3	20.0	27.7	131.3	128.5	136.3	90%
4	13.7	20.3	73.2	70.9	77.7	100%

Table 3: Size  $10 \times 10$ , 4-connectivity,  $N_{it} = 2$ ,  $N_\beta = 3$   $N_{itGibbs} = 6$ 

$N_{lab}$	$E_{DPA}$	$E_G$	$V_{DPA}$	$V_G$	$V_{opt}$	$DPA/Gibbs$
2	14.6	15.9	288.9	292.6	298.0	30%
3	15.2	18.2	155.5	158.1	164.3	20%

Table 4: Size  $8 \times 8$  or  $6 \times 6$ , 8-connectivity,  $N_{it} = 3$ ,  $N_\beta = 4$   $N_{itGibbs} = 12$ 

$N_{hyp}$	$E_{prun}$	$V_{prun}$	$V_{opt}$
256	11.2	128	129.3
16	22.7	123	129.3

Table 5: Size  $8 \times 8$ , 2 labels  $\Rightarrow$  512 running hypotheses.

algorithmic approach to it. In the first place, the experimental monotonicity of the power iteration method that was mentioned in the previous section is somewhat intriguing. Secondly, it is but natural to ask about interior point methods, and one may hope that such methods can be provided by perturbing the polynomial to be maximized so as to make it irreducible.

Finally, in connection with the original motivation presented in section 2, more ambitious questions arise about what happens if the constraint on the degree is relaxed. Deterministic Pseudo-Annealing has been used for a variety of applications, and proven an efficient parallel and deterministic substitute to stochastic methods like Simulated Annealing. While uniqueness of the solution to the deformed problem was stated experimentally a long while ago, the present paper still contributes the establishment of the method by making sure that no pathological situation may arise at this stage. However, when performing subsequent steps, the rate of decreasing of the degree while restoring the original constraints may change the solution we reach, and no theoretical foundations are presently available about such continuation methods.

## References

- [1] M. Berthod. Definition of a consistent labeling as a global extremum. In *International Conference on Pattern Recognition*, pages 339–341, Munich, 1982.
- [2] M. Berthod, S. Liu-Yu, and J.P. Stromboni. Deterministic pseudo-annealing : a new optimization scheme applied to texture segmentation. In *International Conference on Pattern Recognition*, volume 2, pages 533–536, The Hague, Netherlands, sep 1992.
- [3] J. E. Besag. Spatial interaction and the statistical analysis of lattice systems (with discussion). *Journal of Royal Statis. Society B.*, pages 192–236, 1974.
- [4] H. Derin, H. Elliott, R. Cristi, and D. Geman. Bayes smoothing algorithms for segmentation of binary images modeled by markov random fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6, 1984.



- [5] O.D. Faugeras and M. Berthod. Improving consistency and reducing ambiguity in stochastic labeling: an optimization approach. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 4:412–423, 1981.
- [6] F.R. Gantmacher. *The theory of matrices*, volume I,II. Chelsea, 1959.
- [7] S. Geman and D. Geman. Stochastic relaxation, Gibbs distributions and the Bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6:721–741, 1984.
- [8] V. Guillemin and A. Pollack. *Differential topology*. Prentice-Hall, 1974.
- [9] R. Hummel and S. Zucker. On the foundations of relaxation labeling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 5(3), 1983.
- [10] A. Rosenfeld, R.A. Hummel, and S.W. Zucker. Scene labeling by relaxation operations. *IEEE Transactions on Systems, Man and Cybernetics*, 6:420–433, 1976.
- [11] L.G. Shapiro and R. Haralick. Structural description and inexact matching. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 3(5):504–519, 1981.



---

Unité de recherche INRIA Lorraine, Technopôle de Nancy-Brabois, Campus scientifique,  
615 rue du Jardin Botanique, BP 101, 54600 VILLERS LÈS NANCY  
Unité de recherche INRIA Rennes, Irisa, Campus universitaire de Beaulieu, 35042 RENNES Cedex  
Unité de recherche INRIA Rhône-Alpes, 46 avenue Félix Viallet, 38031 GRENOBLE Cedex 1  
Unité de recherche INRIA Rocquencourt, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex  
Unité de recherche INRIA Sophia-Antipolis, 2004 route des Lucioles, BP 93, 06902 SOPHIA-ANTIPOLIS Cedex

---

Éditeur  
INRIA, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex (France)  
ISSN 0249-6399